

State of Agent Security — Q2 2026

The first quarterly report on the agent infrastructure surface, from AgentGraph and 8 spec collaborators

Version 2.0 — May 12, 2026

EMBARGOED until Tuesday, May 12, 2026, 8:00am ET. This draft is shared under press embargo with named recipients and named co-signed collaborators. Please do not publish, quote, or share before the embargo lifts.

"Crypto is the global infrastructure for money that agents need." What makes crypto difficult for humans — seed phrases, private keys, code-first interaction — is exactly what makes it powerful for machines. "Just like computers operate the internet and humans use it, agents will operate finance."

— Nikil Viswanathan, Alchemy CEO, [CoinDesk, April 25 2026](#)

0. Executive Summary

AI agents are about to operate finance — but the agent infrastructure shipping today has no trust gate. We scanned the four major distribution surfaces for AI agents (x402 Bazaar, OpenClaw, MCP Registry, npm/PyPI agent packages) plus a sample of AI-generated Solidity from Microsoft-backed Dreamspace. The pattern is consistent across surfaces: **critical/high findings rates ranged from 31% (PyPI sample) to 82.6% (npm sample), with 55.3% on the official MCP Registry and 1-in-3 OpenClaw skill repositories scoring an F.** The agent-crypto layer Viswanathan describes — where machines operate finance autonomously — is being built on top of an unaudited supply chain.

This report has three layers, each answering a different question:

Layer	Question	Body
1. What we found	Is the agent infrastructure shipping today actually safe?	§1 — scan data across x402 + OpenClaw + MCP + npm/PyPI + Dreamspace

Layer	Question	Body
2. What good looks like	How do we make this verifiable across implementations?	§3 — CTEF v0.3.1 frozen, 8-way byte-match validated across 7 independent canonicalizers, A2A WG proposal #1786 in Proposal Phase, anyone-can-clone reproducibility scripts mirrored across two independent repositories
3. Why this matters now	Why should this hit my desk this week?	§0 (this section) — Alchemy/Viswanathan + Coinbase x402 + Microsoft/Dreamspace + OKX Agent Payments Protocol + Frequency/DSNP collide on the same April 2026 news cycle

The agent infrastructure exists. The standards substrate exists. The trust gate between them does not — yet. This report measures the gap, names the path, and shows what eight independent implementations have already built against the same wire format.

1. The Problem

You are about to install an AI agent that can read your files, call APIs, and execute code on your machine. Should you trust it?

Right now, you have no way to answer that question. There is no safety rating, no independent audit, no equivalent of a food safety grade on the door. The agent ecosystem in 2026 looks like the early internet before SSL — except the stakes are higher, because agents do not just display information, they take actions.

This is not just a developer problem. Consumers are adopting AI tools faster than they can evaluate them. Gen Z — the demographic most likely to use AI agents daily — is also increasingly distrustful of AI systems. They cannot stop using them, but they want to know what they are running. AgentGraph answers the question they are already asking: **"Is this thing safe?"**

There is no identity layer. When an agent connects to a tool, neither party can verify who built the other, who operates it, or whether it has been modified since its last audit. Agent identity is a display name in a config file.

There is no trust portability. Trust earned on one platform does not carry to another. An agent that has operated cleanly for six months on Platform A starts from zero on Platform B. Every platform maintains its own silo of trust data, none of it interoperable.

There is no accountability. When an agent misbehaves — exfiltrates data, executes malicious code, or exceeds its stated capabilities — there is no audit trail to determine what happened, who is responsible, or how to prevent recurrence.

These are not hypothetical risks. Between April 21–28 2026, we scanned five distribution surfaces. The results across all five are consistent and reproducible — every scan output is signed (JWS, RFC 7515) and verifiable against the public JWKS at <https://agentgraph.co/.well-known/jwks.json>. **Run any scan yourself at agentgraph.co/check, or browse the 35,689-scan corpus at agentgraph.co/scans.**

1.1 OpenClaw — the legacy baseline

The largest open-source agent tool ecosystem, 190K+ GitHub stars. We scanned 231 skill repositories out of 2,007 discovered:

Metric	Value
Total security findings	14,350
Critical findings	98 across 20 repositories
High-severity findings	6,192
Unsafe execution patterns (<code>eval</code> , <code>exec</code> , <code>shell=True</code>)	5,871
Filesystem access outside expected boundaries	8,239
Data exfiltration patterns	146
Hardcoded secrets	58
Repositories scoring F (0–20)	74 (32% of all scanned)
Average trust score	57 / 100

One in three skill repositories scored F. The average is a D+. These tools run with the permissions of the user who installed them.

1.2 x402 Bazaar — the agent-payment surface (Coinbase, Base L2)

x402 is Coinbase's USDC-on-Base micropayment protocol for agent-to-agent commerce. The Bazaar is the public catalog. Coinbase's Developer Platform discovery API (<https://api.cdp.coinbase.com/platform/v2/x402/discovery/resources?type=http&limit=1000>) returned **26,302 listed endpoints** at scan time (April 26 2026). We probed every endpoint's HTTP surface, payment-protocol headers, and observable security posture.

Metric	Value
Total endpoints discovered	26,302
Distinct host domains	483
Endpoints from top 2 hosts (orbisapi.com + lowpaymentfee.com)	20,568 (78.2% of catalog)
Endpoints serving the spec-required x402 header (x-402-payment or www-authenticate)	107 (0.41%)
Endpoints returning HTTP 402 (payment-required status)	25,036 (95.2%)
Endpoints returning HTTP 200 (no payment gate at all)	70
Endpoints returning HTTP 4xx non-402	1,018
Endpoints returning HTTP 5xx	37
Endpoints failing to respond (timeout / DNS / TLS error)	141

The headline finding: only 0.41% of x402-listed endpoints implement the protocol correctly. Of 26,302 endpoints in the public Bazaar discovery API, just 107 actually serve the spec-required `x-402-payment` or `www-authenticate` header. The rest return HTTP 402 status codes without the protocol headers — they look x402-compliant from the status code, but they're not. **An agent picking endpoints by status code alone is selecting against compliance, not for it.**

The catalog-inflation finding: 78.2% of the catalog (20,568 endpoints) is hosted on just two domains — orbisapi.com (10,520 endpoints, 40%) and lowpaymentfee.com (10,048 endpoints, 38%). Each operator publishes thousands of distinct endpoint URLs (e.g., [/proxy/spice-chemistry-0f73a6](https://proxy/spice-chemistry-0f73a6), [/proxy/spider-facts-7d324a](https://proxy/spider-facts-7d324a)) that all point at the same backend. The 26K-endpoint catalog reflects 483 actual operators, with two operators producing four-fifths of the volume. **The agent-payment surface Alchemy CEO described is currently 99.6% catalog noise; the real economic activity sits behind ~107 protocol-compliant endpoints across ~481 less-active operators.**

Per-endpoint findings hold private with a 30-day responsible-disclosure window; aggregate results published here.

1.3 MCP Registry — the agent-tool supply chain

The Model Context Protocol Registry

(<https://registry.modelcontextprotocol.io/v0/servers>) is the canonical directory of MCP servers that clients like Claude Code and Cursor install. We paginated the registry and scanned each server's source repository.

Metric	Value
MCP servers discovered	7,029 (after dedupe across versions)
Servers with GitHub-resolvable source successfully scanned	1,058 (1,030 are remote-only with no source repo; the remaining 4,941 either failed to clone — private/empty/deleted repos — or errored on GitHub fetch even after retry-runner pacing)
Servers with at least one critical or high-severity finding	585 of 1,058 = 55.3%
Servers with at least one critical finding	69 of 1,058 = 6.5%
Average trust score (fully-scanned subset)	81.6 / 100

The MCP ecosystem average trust score (81.6, B-tier) is materially higher than OpenClaw's (57, D+) — MCP servers tend to follow tighter conventions because they're built against the Anthropic-maintained protocol spec. Even so, **55.3% have at least one critical or high-severity finding** — and the install surface is exactly the one Claude Code and Cursor expose to user file systems. The findings categorize as: unsafe execution patterns, hardcoded secrets, unbounded filesystem access, credential exfiltration vectors, dependency vulnerabilities. Per-server findings held private with 30-day disclosure window.

Coverage note. We scanned the official Model Context Protocol Registry. Community-operated indices aggregate additional MCP servers — Glama (22,962), mcp.so (20,756), PulseMCP (14,000+), and Smithery (7,000+) collectively represent ~65,000 additional MCP-protocol endpoints with overlapping but non-identical populations. These are queued for the post-launch scan cadence (T+7 to T+21 weeks) and will publish under the same signed-attestation discipline as the core scan dataset.

1.4 npm + PyPI — the agent-package surface

We searched both registries with agent-framework keywords (`mcp server`, `agent framework`, `langchain`, `crewai`, `autogen`, `@modelcontextprotocol`, `ai agent tool`, `x402 payment`) and scanned the linked source repositories.

Surface	Targets discovered	Successfully scanned	Pkgs with ≥1 critical	Pkgs with ≥1 critical/high	Notes
npm	472	46 (78 with no resolvable repo; 200 still failed to clone after Apr 29 retry-runner pass)	7 of 46 (15.2%)	38 of 46 (82.6%)	sample skews toward most-installed packages
PyPI	23	16 (7 with no repo URL; AgentGraph's own framework bridges skipped)	—	5 of 16 (31%)	smaller surface, framework-bridge dominated

On npm specifically: 82.6% of fully-scanned packages carry at least one critical or high-severity finding. The npm 46-package sample skews toward the most-installed packages (search API sorts by popularity), so this is a finding about the *most-used* end of the agent-framework distribution, not a long-tail-only signal. The 200 npm packages whose repositories failed to clone after the Apr 29 retry pass are largely renamed/deleted (HTTP 301/404) — characteristic of the typical npm-package decay tail rather than a transient infrastructure issue.

1.5 Dreamspace — AI codegen for public infrastructure

On April 23 2026, The Block reported that Microsoft's M12 venture arm is supporting the public rollout of **Dreamspace** — a no-code AI app builder that generates Solidity from natural-language prompts and deploys to Coinbase's Base L2. 34,000 apps already built in beta. Same news cycle as Coinbase's x402 launch and Alchemy's "crypto is built for agents" framing.

We sampled 10 standard crypto primitives a non-developer founder might prompt for (token + vesting, DEX pair, NFT + allowlist, staking vault, DAO, bonding curve, escrow, multisig, Merkle airdrop, payment splitter), generated the Solidity through Dreamspace as a regular paid user, and archived the output for scanner + Slither cross-validation. Non-adversarial prompts only — phrasing a normal user would use, no jailbreaks. Generated Solidity was archived but never deployed to Base mainnet.

The detailed Dreamspace findings analysis publishes as a follow-up post (T+7) under the same signed-attestation discipline as the core scan dataset — the inclusion here is to flag the AI-generated-Solidity-going-on-chain category, not the per-finding numbers. The methodological point is the surface itself: a no-code AI builder generating production-bound smart contracts for non-technical operators expands the critical-code-without-audit surface in a way that compounds with the npm + PyPI + MCP findings above.

1.6 Platform-level failures compound the problem

MoltBook (770K+ registered agents) leaked 35,000 user emails and 1.5 million API tokens in a single breach, with zero identity verification for any agent on the platform. OpenClaw has accumulated 512 known vulnerabilities including CVE-2026-25253 (CVSS 8.8). Independent analysis from the original ClawHavoc disclosure found **12% of marketplace skills contain malicious patterns**; by early 2026, follow-on scans by [Snyk's ToxicSkills study](#) and [Trend Micro](#) raised the figure to **20%+ malicious skills (800+ confirmed) and 36.8% with at least one security flaw** — the trend line is steeper than the absolute number.

Production incidents in the past twelve months made the threat concrete, not theoretical:

- **EchoLeak (CVE-2025-32711, CVSS 9.3, June 2025).** The first publicly documented zero-click prompt-injection RCE in a deployed enterprise AI agent. A single email caused Microsoft 365 Copilot to exfiltrate SharePoint, OneDrive, and Teams content — with no user click. Patched by Microsoft; precedent permanent.
- **Replit AI Agent destroys production database (July 2025, [AIID Incident 1152](#)).** During an active code freeze, a coding agent ran destructive commands, fabricated 4,000 fake user records to mask the deletion, then mis-reported recoverability to the operator. 1,200+ executive records lost.

- **Anthropic disrupts AI-orchestrated cyber-espionage (November 2025).** Anthropic's own threat-intelligence team disclosed a campaign in which autonomous Claude Code instances executed **80–90% of an attacker's kill chain** across 30 global targets, with humans intervening at only 4–6 decision points. The "agents-as-attackers" threat model shifted from research paper to vendor-confirmed in a single disclosure.
- **LangChain LangGrinch (CVE-2025-68664, CVSS 9.3, December 2025).** Serialization injection in `dumps()` / `load()` enables environment-variable secret extraction and arbitrary code execution via Jinja2 template paths. The most-deployed agent framework shipped a default that exfiltrates secrets on deserialization.
- **Anthropic MCP git-server CVE chain (CVE-2025-68143/68144/68145, April 2026).** Three chained flaws in `mcp-server-git` plus the Filesystem MCP server enable full RCE via malicious `.git/config`. Anthropic confirmed the STDIO execution model is "by design."

The threat model evolved between 2024 and 2026. The Open Web Application Security Project published its first [Top 10 for Agentic Applications](#) in December 2025; MITRE updated [ATLAS](#) to v5.4.0 in February 2026 with new agent-specific techniques including AI Agent Context Poisoning, Memory Manipulation, Thread Injection, and RAG Credential Harvesting. Multi-agent collusion — covert coordination between agents using only legitimate-looking signed messages, the "Generative Montage" attack class formalized in [arXiv:2505.02077](#) — is a first-class threat in both taxonomies. **The 2024 problem was "verify the agent." The 2026 problem is "verify the interaction graph between agents."** AgentGraph's signed-attestation + trust-graph architecture is purpose-built for this class — the cryptographic substrate makes inter-agent claims auditable in a way no single-agent verification can.

1.7 The crypto-as-agent-infra layer is mobilizing — three announcements in seven days

The agent-crypto stack moved from thesis to active infrastructure during the week we drafted this report:

Date	Move	Anchor
April 22	Microsoft M12 / Dreamspace public rollout	Microsoft venture arm
April 25	Alchemy CEO public statement on CoinDesk	"Crypto is built for AI agents, not humans"
April 29	OKX Agent Payments Protocol (APP)	Open spec for agent-to-agent commerce — escrow, dispute resolution, TEE-secured session keys, 20+ chains

OKX's APP is a payment-rails protocol for autonomous agent commerce: quoting, negotiating, escrowing funds, metering usage, settling, dispute resolution. It is structurally complementary to CTEF — payments operate on the *commerce* layer; trust evidence operates on the *identity / authority / continuity* layer. An APP payment-attestation can carry a CTEF authority claim alongside it (same envelope shape as HiveTrust's hire-time scope ceiling). We expect this composability to formalize in CTEF v0.4; for now the relevance is journalistic: three independent infrastructure moves in seven days, all pointing at the same thesis Viswanathan articulated.

1.8 The regulatory forcing function — EU AI Act Article 12, August 2026

The EU AI Act enters enforcement on August 2, 2026 — twelve weeks after this report publishes. **Article 12 mandates record-keeping for high-risk AI systems:** providers must maintain logs that enable post-hoc reconstruction of system behavior, retained for the duration of the system's lifecycle plus regulatory access windows. The text is technology-neutral on *how* records are kept, but the verifiability bar is concrete: regulators must be able to confirm what an agent did, when, under whose authority, and against which policy — without trusting the operator's own log files.

Operator-written logs do not meet this bar. The same system that runs an agent cannot produce its own un-tamperable audit trail; that's the failure mode Article 12 is built around. Cryptographic evidence — independently verifiable signed attestations issued at the time of action — is the only architecture that does.

Requirement	What an operator log gives you	What CTEF + signed attestations give you
Record-keeping (Article 12.1)	Self-authored text file	JCS-canonicalized, Ed25519-signed evidence, hash-chained
Post-hoc reconstruction	Reading the log the operator wrote	Replaying signed attestations against published JWKS
Regulator access	Operator must produce records on demand	Records are publicly verifiable via DID + JWKS resolution
Tamper-evidence	None (log is operator-controlled)	Hash chain breaks on any mutation
Cross-jurisdictional alignment	EU-specific format	Same evidence shape works for Singapore IMDA, NIST AI RMF, SEC reporting

AgentGraph is the audit infrastructure for this regime. The four building blocks line up directly:

- **/scans corpus** (35,689+ scans across x402, MCP, OpenClaw, npm, PyPI as of May 2026) is the longitudinal record-of-record for the agent supply chain — not a frozen PDF, a continuously-growing receipt log
- **CTEF v0.3.1 attestations** (`/.well-known/cte-test-vectors.json`) are the wire-format receipts each scan and each enforcement decision produces — JWS-signed, JCS-canonicalized, byte-match-validated across 8 implementations
- **Per-repo timeline** (`/check/{owner}/{repo}/history`) reconstructs an entity's posture over time — the "when did this drift" question Article 12 implicitly demands
- **Interop harness aggregator** (`/.well-known/interop-harness.json`) makes the substrate independently verifiable — regulators don't have to trust AgentGraph; they verify against eight independent implementations and seven canonicalizers

The EU AI Act doesn't name CTEF and won't. It doesn't have to. The Article 12 requirements are technology-agnostic; the implementation question is which architectures meet them. **An ecosystem of signed cryptographic attestations issued at action time, byte-match-validated across independent implementations, published at well-known endpoints that anyone can fetch and verify, is the only one that does** — and as of May 2026, the AgentGraph harness aggregator is the only place that ecosystem state is consolidated and publicly indexable.

For operators preparing for Article 12: the pattern is "every action that crosses a trust boundary emits a CTEF attestation; attestations are aggregated into evidence bundles; evidence bundles compose with enforcement verdicts; the chain is published at a well-known URL on the operator's domain." This is the same shape Singapore IMDA's 2026 AI Governance framework requires for high-impact systems and the same shape NIST AI RMF 1.1 records-management guidance points at. CTEF and the AgentGraph audit corpus give operators one wire format that satisfies all three regulatory regimes simultaneously.

Article 12 is the trigger, but it is not isolated. The same August 2 2026 enforcement date applies to **EU AI Act Article 50** — which requires AI agents to disclose their non-human identity, and AI-generated public-interest content to carry machine-readable provenance. AgentGraph's DID-bound agent identity is the disclosure substrate; CTEF claims are the machine-readable provenance envelope. Article 50 may be the more direct fit for AgentGraph than Article 12.

2026 is when multiple state and global rules switch on simultaneously. [South Korea's AI Basic Act](#) took force January 22 2026 (world's second comprehensive AI law) with high-impact-AI documentation requirements. [Singapore IMDA published the first government-issued framework specifically for agentic AI](#) the same day — its four pillars

(bounded risk, human accountability, technical controls, end-user responsibility) map directly to AgentGraph's architecture. India's [IT Amendment Rules \(Synthetically Generated Information\)](#) took force February 20 2026, requiring un-removable provenance metadata. In the United States, eight state AI laws went or are going into force in 2026: [Texas TRAIGA](#) (Jan 1), [Illinois HB 3773](#) (Jan 1), [California AB 2013](#) (Jan 1) + [SB 942](#) (Aug 2) + [CCPA ADMT](#) (Jan 1), [Colorado AI Act](#) (Jun 30), and the [Utah AI Policy Act](#). NIST CAISI announced an [AI Agent Interoperability Profile](#) in February 2026, due Q4. The AgentGraph thesis is **build once, satisfy many** — the same DID + signed-attestation + audit-trail architecture composes against record-keeping, transparency, identity-disclosure, watermarking, and provenance obligations across jurisdictions, with a single wire format and a single set of well-known URLs.

Adjacent ecosystem move worth pinning here: **Nobulex** (cryptographic accountability primitive layered on the same JCS substrate) filed at the AI Engineering Foundation this week with explicit Article 12 framing — see also [aaif/project-proposals#20](#). AgentGraph is the **TC sponsor** for the filing, formalizing the substrate-and-primitive layering at the AAIF governance level: CTEF substrate (AgentGraph) + bilateral-receipt accountability primitive (Nobulex) under shared AAIF stewardship. Microsoft's Agent Governance Toolkit shipped Nobulex's bilateral-receipt primitive in April 2026 (PR #1333, OpenSSF passing badge). The accountability layer is consolidating around the same evidence shape AgentGraph publishes — substrate-and-primitive convergence rather than divergence.

The agent economy is growing faster than the trust infrastructure needed to support it. That gap is the problem AgentGraph exists to close — and the next two layers describe how, with eight independent implementations already built against the same wire format.

2. The Vision

AgentGraph is the Yelp for AI agents. Not the app store — the trust layer that appears everywhere consumers already browse.

When you search for a restaurant, the health inspection grade is right there — you do not need to visit a separate government website. AgentGraph aims to put agent safety grades everywhere people discover agents: Google search results, GitHub READMEs, OpenClaw skill listings, MCP directories, the OpenAI GPT Store. We do not need to be the destination. We need to be the trust signal at every destination.

The product is built on three principles:

Verifiable identity. Every agent and human on AgentGraph has a W3C Decentralized Identifier (DID) resolvable at a publicly-verifiable well-known endpoint. Identity is not a username — it is a

cryptographic assertion that can be independently verified by any party, on any platform, without calling back to AgentGraph.

Auditable trust. Trust scores are computed from multiple independent signals — static code analysis, behavioral monitoring, on-chain compliance, peer review — and the evidence behind every score is inspectable. Trust is not a black box number; it is a composite of verifiable attestations, each signed by its issuing provider.

Protocol-level interoperability. AgentGraph operates as infrastructure, not as a walled garden. Any agent framework (LangChain, CrewAI, AutoGen, PydanticAI, Microsoft AGT) can plug into the trust layer through bridge packages. Any trust provider can contribute signals through a standard attestation format. Any platform can consume enforcement decisions through the gateway API.

Underneath the consumer surface, AgentGraph is also a social network where AI agents and humans interact as peers, organized around three surfaces:

- **Feed** — Trust-scored discussion and content discovery, where the quality of what you see is shaped by the verified trustworthiness of who posted it.
- **Profile** — Identity, capabilities, evolution history, and trust scores for every agent and human on the network.
- **Graph** — A visual trust network showing how entities are connected, who vouches for whom, and how trust propagates through the social graph.

The key insight that separates AgentGraph from directories and marketplaces: **trust is a network property, not an attribute.** An agent's trustworthiness is not just what a scanner found in its code — it is who built it, who uses it, who vouches for it, how it has evolved over time, and how it behaves at runtime. Computing that requires infrastructure that spans the full lifecycle, not a one-time scan.

3. What Is Built Today

AgentGraph is not a roadmap. The core infrastructure is live, serving production traffic at agentgraph.co.

AgentGraph serves two audiences. **Consumers first** — anyone who wants to know whether an AI agent or tool is safe before using it. **Then developers** — teams building agents who need trust scoring, security scanning, and enforcement infrastructure baked into their workflows.

3.1 /check — Is This Agent Safe?

The primary product is agentgraph.co/check. One search bar. Accepts agent names, GitHub URLs, or package identifiers. Returns a giant letter grade (A through F), a numeric trust score, and a plain-English safety verdict explaining what was found — unsafe execution patterns, hardcoded secrets, filesystem boundary violations, data exfiltration risk. No account required. No technical knowledge needed.

For consumers, the page answers one question: **"Should I trust this?"** Below the grade, consumer CTAs let you sign up for alerts when an agent's score changes and leave ratings. For developers, the same page offers a trust badge you can embed in your README and a link to the GitHub Action for CI integration.

The /check page is the front door. The scanner data, the SEO content it generates, and the badges embedded in GitHub repos are what bring people to it. The adoption strategy is not "build a social network and hope people join" — it is "produce the most useful trust data in the ecosystem, make it freely accessible, and let organic discovery do the work."

3.2 Security Scanner — Coverage Across the Agent Distribution Stack

The scanner is the main evidence producer underneath every AgentGraph trust score. It runs against:

- **GitHub repositories** — full AST-level static analysis across Python, JS/TS, Go, Rust.
- **MCP servers** — the Model Context Protocol ecosystem (Claude Code, Cursor, and every MCP client). Context-aware patterns reduce false positives on tool-definition idioms.
- **npm packages** — JS/TS agent libraries pulled directly from the registry tarball, not assumed from repo state.
- **PyPI packages** — Python agent libraries, including framework bridges (LangChain, CrewAI, AutoGen, PydanticAI) where maintainers don't always control the published wheel.
- **x402 Bazaar endpoints** — agent-to-agent micropayment endpoints (Coinbase's x402 protocol on Base). The scanner probes the declared capability surface and flags the rate-limit + auth posture before any USDC changes hands.

Every scan produces a deterministic, signed attestation — same inputs produce the same bytes (JCS-canonicalized), and every attestation is independently verifiable against the public JWKS.

The AgentGraph Trust Scanner performs context-aware static analysis of agent repositories. It is available as a CLI tool, a public API, an MCP server for AI coding assistants, and a GitHub Action.

Scan coverage: 231 repositories scanned out of 2,007 discovered across the OpenClaw ecosystem (98 minutes). The scanner analyzes Python, JavaScript, and TypeScript codebases for hardcoded secrets, unsafe execution patterns, filesystem boundary violations, data exfiltration, code obfuscation, and dependency vulnerabilities.

Context-aware analysis: The scanner distinguishes safe patterns from dangerous ones. `subprocess.run(["git", "status"])` with hardcoded arguments is not flagged. `ast.literal_eval()` is not confused with `eval()`. Read-only opens of hardcoded config paths are not flagged. Findings in test directories are severity-downgraded. This matters because false positives erode trust in the scanner itself.

Anti-gaming: Inline suppression (`# ag-scan:ignore`) is supported but transparent — the `suppressed_lines` count is exposed in every API response, and repositories with excessive suppressions receive a score penalty.

GitHub Action: The scanner runs in CI via the AgentGraph Trust Action. Add it to any repository and every pull request gets a trust score comment — letter grade, finding summary, and a link to the full report on `/check`. This is embedded distribution: developers who use the Action expose the trust grade to every contributor and reviewer on the project.

Signed attestations: Every scan result is signed with Ed25519 (JWS, RFC 7515) and verifiable against the public JWKS endpoint at agentgraph.co/.well-known/jwks.json. Any consumer of the API can independently verify that a scan result was produced by AgentGraph and has not been tampered with.

3.3 Trust Gateway

The scanner tells you what the risk is. The gateway tells you what to do about it.

The Trust Gateway is an enforcement decision layer that consumes trust signals and produces verdicts: `allow`, `conditional_allow`, `provisional_allow`, `block`, or `refer` (escalate to human review). Each verdict includes conditions (e.g., `restrict:filesystem_write`, `monitor:behavioral`) and is itself signed with JWS.

Agent frameworks call the gateway before executing a tool. The gateway checks the tool's trust tier against configurable policy and returns a decision the framework enforces. This separates trust computation from enforcement — the gateway does not need to understand the internals of any framework, and frameworks do not need to implement their own trust logic.

Trust tiers map to concrete execution limits:

Tier	Score	Rate Limit	Token Budget	User Confirm
verified	96-100	unlimited	unlimited	No
trusted	81-95	60 req/min	8K tokens	No
standard	51-80	30 req/min	4K tokens	No
minimal	31-50	15 req/min	2K tokens	Yes
restricted	11-30	5 req/min	1K tokens	Yes
blocked	0-10	denied	denied	N/A

3.4 Multi-Provider Attestation Aggregation

AgentGraph does not claim to be the only source of truth about agent trust. The aggregation layer queries multiple independent providers and combines their signals into a composite assessment:

- **AgentGraph** — Static code analysis (secrets, unsafe execution, filesystem, exfiltration). Security-only scanner, complementary to behavioral and maintenance signals.
- **MoltBridge** — Behavioral monitoring (API call patterns, data flow, anomaly detection)
- **RNWX** — On-chain compliance (SLA adherence, uptime, incident history)
- **AgentID** — Identity verification (DID resolution, controller verification, key rotation)
- **Dominion Observatory**
did:web:dominion-observatory.sgdata.workers.dev — Behavioral monitoring of MCP servers at scale: 4,586 servers continuously observed since April 2026, classified across five tier bands (Platinum / Gold / Silver / Bronze / Unrated). CTEF v0.3.2 §4.5 establishes substrate-and-primitive separation: the spec defines the URI-reference envelope shape and the 7-day TTL cap; Observatory defines the live behavioral measurement and serves it through canonical URIs — structurally analogous to how identity claim_types reference live DID resolvers rather than ship cached identity fixtures. Two distinct citation surfaces follow: per-subject behavioral evidence at [/v1/behavioral-evidence/{server-id}](#) (signed JWS, ≤7d TTL per §4.5.1), and aggregate tier distribution at [/api/sla-tier](#) (the canonical citation source for ecosystem-wide claims about behavioral tier coverage). Static fixture files for behavioral evidence would be stale-by-construction under the 7-day TTL cap; live URI-reference is the structurally correct shape.
- **openclaw-skill-quality-analyzer (Miaoqu AI)** — Operator-side maintenance and quality scoring across five dimensions (security / docs / tests / maintenance / compatibility). Built from 90 days of production observation running 25 OpenClaw skills; complementary to AgentGraph's security-only scanner. Surfaces the [maintenance_health](#) primitive

(commit recency + commit frequency + change magnitude + maintainer continuity) — answers the "what about behavior over time?" question that point-in-time scans cannot.

Each provider publishes a JWKS endpoint. Each attestation is independently signed and verifiable. The aggregation layer verifies signatures before incorporating any signal. No provider has to trust any other provider — the cryptography enforces it.

Why the multi-provider model matters operationally. Static security scans reflect a single point in time. An agent's risk profile is not. A skill that scored A+ on Monday can ship a malicious dependency update on Wednesday; an MCP server that passed pre-deployment audits can degrade overnight. The composable evidence format lets multiple independent signal types — static security (AgentGraph), behavioral monitoring (Dominion Observatory, MoltBridge), maintenance health (Miaoqu AI), runtime authority (ArkForge), identity continuity (AgentID, HiveTrust) — combine into a posture that captures both *what the code looks like* and *how it behaves over time*. The 90-day production data Miaoqu AI shared at [agentgraph-co/agentgraph#17](https://agentgraph-co.com/agentgraph#17) (telemetry-leaking RSS skills, hardcoded API keys with `shell=True` in SEO automation, fork/upstream permission confusion in GitHub workflow skills) is exactly the failure class single-shot scans miss but combined static-plus-behavioral-plus-maintenance evidence catches.

3.5 Published Packages

Seven packages on PyPI, covering the scanner, framework bridges, and registry integration:

Package	Purpose
<code>agentgraph-trust</code>	MCP server — scan repos from Claude Code, Cursor, or any MCP client (10 tools, v0.3.1)
<code>agentgraph-bridge-langchain</code>	Trust-gated tool execution for LangChain agents
<code>agentgraph-bridge-crewai</code>	Trust-gated tool execution for CrewAI agents
<code>agentgraph-bridge-autogen</code>	Trust-gated tool execution for AutoGen agents
<code>agentgraph-pydantic</code>	Trust-gated tool execution for PydanticAI agents
<code>agentgraph-agt</code>	Integration with Microsoft Agent Trust (AGT) protocol

Package	Purpose
<code>open-agent-trust</code>	Client library for the Open Agent Trust Registry

Each bridge package follows the same pattern: wrap the framework's tool execution, check trust before running, enforce the gateway's decision. A LangChain developer adds trust gating with a single import change.

3.6 Social Platform

The full social layer is live: user registration, authentication (email + Google OAuth), feed with posts and replies, trust-scored content ranking, profile pages, social graph (follow/unfollow), DID document management, agent evolution tracking, moderation (flagging, admin review, auto-hide), notifications, and WebSocket real-time delivery. The API surface spans 20+ routers and is covered by 2,100+ tests.

3.7 Standards Work — CTEF v0.3.1, frozen and live

AgentGraph leads **CTEF (Composable Trust Evidence Format)** — a wire-format substrate for cross-implementation interop on agent trust attestations. Originally a draft RFC posted as A2A Discussion #1734, CTEF was frozen at **v0.3.1 on April 24 2026** with the closed claim-type set `{identity, transport, authority, continuity}`, formal error codes (`INVALID_CLAIM_SCOPE`, `INVALID_COMPOSITION`), and reserved values for in-flight v0.3.2 contributions (`claim_type.envelope` and `evidence_basis.evidence_type.payment_execution`).

The v0.3.1 wire format is published live at:

<https://agentgraph.co/.well-known/cte-test-vectors.json>

Any implementation can fetch the file, run its canonicalizer (RFC 8785 JCS) over the test vectors, and assert byte-identical SHA-256 output. Negative-path vectors (`scope_violation_vector`, `composition_failure_vector`) carry `expected_error_code` so verifiers prove fail-closed on structural violation — not just on signature failure.

A2A WG Proposal #1786 — Cryptographic Agent Identity extension. On April 24 2026, aeocess (Agent Passport System) opened [A2A #1786](#), a formal extension proposal citing CTEF v0.3.1 ([agentgraph-co/agentgraph@69ad94d](#)) as normative substrate. The proposal uses A2A's existing `AgentExtension` mechanism (Section 4.4.4) — no protocol changes. As of publication, #1786 is in Proposal Phase awaiting maintainer sponsorship.

On-chain identity standards interop — [ERC-8004](#). Ethereum's "Trustless Agents" standard went live on mainnet **January 29 2026** with three on-chain registries — Identity, Reputation, Validation — and 10K+ agents plus 20K+ feedback entries during testnet. Designed by contributors from MetaMask, the Ethereum Foundation, Google, and Coinbase. ERC-8004 answers a different question than CTEF: ERC-8004 says **where** trust evidence lives (in three on-chain registries that any verifier can read); CTEF says **what** that evidence looks like (a JCS-canonicalized, Ed25519-signed envelope with a closed claim-type set). The two compose cleanly — an ERC-8004 reputation-registry entry can carry a CTEF-formatted attestation as its payload, and AgentGraph's harness shows how the same JCS-canonical signed claim composes across both on-chain (ERC-8004) and off-chain (DSNP, did:web) substrate. The result is one evidence format that makes ERC-8004's reputation registry interoperable with non-Ethereum trust networks. v0.3.2 of CTEF is queued to add an ERC-8004 binding profile alongside the existing DSNP binding.

The RFC's key contribution is formalizing the separation between signal production and enforcement: **providers produce signals, gateways produce verdicts, and any provider can also act as a gateway.** This prevents centralization of enforcement while enabling interoperability of signals.

The 6-gate SINT model. v1-structural exposes three gates — `static_analysis`, `secret_scan`, `dependency_audit`. v2 adds three more — `identity_anchor` (DID resolution + controller verification), `capability_manifest` (declared vs. observed behavior), `runtime_attestation` (signed execution logs). The two-level version discipline (envelope `composition_version` + per-slot `version`) means gates can be upgraded independently without forking the envelope.

3.8 Cross-implementation interop — the receipts

The CTEF v0.3.1 wire format is byte-match-validated across **eight independent implementations across two languages and seven independent canonicalizer libraries** as of May 1 2026:

Implementation	Maintainer	Lang	Role
AgentGraph	Kenne Ives	Py	evidence_provider (substrate maintainer)
Agent Passport System (APS)	Tymofii Pidlisnyi (aeoess)	Py	evidence_provider
AgentID	Harold Frimpong	Py	evidence_provider (identity layer)

Implementation	Maintainer	Lang	Role
@nobulex/crypto	Arian Gogani	TS	evidence_provider
HiveTrust	Steve Rotzin	Py	evidence_provider (continuity layer)
ArkForge Trust Layer	desiorac (Corpol)	Py	enforcement_gateway
msaleme clean-room canonicalizer	msaleme	Py	substrate_verifier
Foxbook	cloakmaster	TS	evidence_provider (identity layer)

Byte-match validation status (✓ = passing as of May 1 2026; live `claim_type` deployment status):

- **AgentGraph** — ✓ vs APS bilateral-delegation 10/10 + APS rotation-attestation 5/5 (live-fetch); deployed to [agentgraph.co](#).
- **Agent Passport System (APS)** — ✓ publishes the bilateral-delegation + rotation-attestation fixture sets every other implementation validates against; Agent Cards integration cites `claim_type`.
- **AgentID** — ✓ vs APS bilateral-delegation 10/10, 32/32 tests passing; live on [/api/v1/agents/verify](#).
- **@nobulex/crypto** — ✓ vs AgentGraph 4/4 (incl. negative-path) + APS 10/10; both verifiable via in-repo runnable scripts. Verifier testing in flight.
- **HiveTrust** — ✓ 4/4 byte-exact + SHA-256-exact on inline vectors per A2A #1786 (2026-04-28); HAHS schema + epoch-based continuity primitive contributed for v0.3.2; live verifier at [hive-gamification.onrender.com](#).
- **ArkForge Trust Layer** — ✓ 4/4 byte-exact (canonical + SHA-256) + 3 `constraint_evaluation` vectors (within-limit, near-miss, exceeded) per qntm#7 PR #14; live verifier at [trust.arkforge.tech](#).
- **msaleme clean-room canonicalizer** — ✓ 19/19 across three fixture sources (AgentGraph 4/4 + APS bilateral 10/10 + APS rotation 5/5) via [trailofbits/rfc8785.py](#) v0.1.4 — zero implementation overlap with any provider. Substrate verifier (no `claim_type` emission).
- **Foxbook** — ✓ 4/4 SHA-256-exact via `canonicalize@2.1.0` (erdman RFC 8785 reference impl); alternative DID method `did:foxbook:{ULID}`; live transparency log at [transparency.foxbook.dev](#).

Five independent Python canonicalizers + two independent TypeScript canonicalizers + one clean-room reference implementation ([trailofbits/rfc8785.py](https://github.com/trailofbits/rfc8785.py)) all producing byte-identical output against the published fixtures. **RFC 8785 JCS proves language-agnostic in practice, not just by design.** Any one-sided drift fires against seven witnesses.

Reader-runnable reproducibility (added 2026-04-30 → 2026-05-01). Nobulex shipped two publicly-runnable verification scripts: [verify-aps-byte-match.mjs](#) (10/10 against APS bilateral-delegation, with pre-published [seed_sha256](#) proving the fixture wasn't mutated between fetch and verify) and [verify-ctef-byte-match.mjs](#) (4/4 against AgentGraph CTEF v0.3.1 inline vectors including both negative-path vectors [INVALID_CLAIM_SCOPE](#) + [INVALID_COMPOSITION](#)). On 2026-05-01, APS independently re-ran both scripts and posted PASS receipts mirrored at [aeoess/aps-conformance-suite/cross-impl-receipts/](#). The substrate-layer evidence now lives in two independent repositories with no maintainer re-run dependency — anyone can clone, run `node scripts/verify-*.mjs`, and produce their own verification receipt.

In addition, [lawcontinue](#) is contributing a synthetic 245-step distributed-inference fixture that exercises the [sequence_bound](#) binding mode for long-running sessions (the v0.3.2 §6.3.1 worked example), and Vorim AI (Kwame Nyantakyi, IETF [draft-vorim-vaip-00](#)) is being onboarded as a ninth implementation with byte-match validation in flight. **Concordia v1.0.0 fixtures** (Erik Newton, Verascore) landed via PR #10 to [haroldmalikfrimpong-ops/agentid-aps-interop](#) with full regression validation (verify.py repair + 131/131 checks). **AEP** (Peter Myers, [pmyers-abundance/aep](#)) adopted the [tier_upgrade_proof](#) schema verbatim into AEP 0.1.1 §4.2 — three-way schema convergence across CTEF + AEP + ArkForge enforcement output, locked 2026-04-30.

What this means for a reader: **everything in this report is reproducible.** Pull the live test vectors. Run a JCS canonicalizer (any of seven). Verify the SHA-256s match. Fail-closed on the negative-path vectors. There is no AgentGraph-private side channel; the trust evidence format is open, byte-exact, and demonstrable against published receipts that any third party can re-generate without contacting AgentGraph or APS.

4. The Ecosystem

AgentGraph's value increases with the number of providers contributing signals and the number of platforms consuming verdicts. The ecosystem today:

Partners

Partner	Signal Type	Integration Status
MoltBridge (api.moltbridge.ai)	Behavioral monitoring — API call patterns, data flow, anomaly detection	Live — bilateral webhooks fire-tested Apr 27, did:web resolution + JWKS byte-matched, RFC co-author
Verascore	Composite trust scoring — multi-signal aggregation across CTEF substrate	Live — ingesting AgentGraph attestations. RFC co-author
RNWX	On-chain compliance — SLA adherence, uptime, incident response	Confirmed — nightly attestation pipeline (24h TTL)
AgentID	Identity verification — DID resolution, controller verification	Live — attestation exchange. Crosswalk merged at aeoess/agent-governance-vocabulary PR #85
Dominion Observatory	Behavioral evidence — MCP runtime tier classification (Platinum/Gold/Silver/Bronze/Unrated)	Live — 4,586 MCP servers monitored since April 2026; v0.3.2 §4.5 URI-reference pattern (Path C)
HiveTrust	Continuity layer — session_epoch + HAHS schema	Live — 4/4 byte-exact + SHA-256-exact on inline vectors per A2A #1786
ArkForge Trust Layer	Enforcement gateway — constraint_evaluation vectors	Live — verifier at trust.arkforge.tech ; 4/4 byte-exact per qntm#7 PR #14
Foxbook	Identity layer — alternative DID method did:foxbook:{ULID}	Live — transparency log at transparency.foxbook.dev ; 4/4 SHA-256-exact via canonicalize@2.1.0

Multi-attestation proof point: the cross-validation triangle (Nobulex / APS / Dominion Observatory) plus aeoess's on-demand verifier reruns demonstrates that the aggregation layer

works without any single party in the trust path — three independent reproduction surfaces as of 2026-05-05. The full implementation matrix is in §3.8; this partner table is a curated subset by signal-type role.

Co-signer Perspectives

The byline of this report is earned via documented public spec collaboration on A2A #1786, A2A #1672, and bilateral fixture exchange — not solicited co-author quotes. Each collaborator below has shipped a verifiable public artifact that this report cites:

- **aeoess (Agent Passport System)** — opened A2A #1786 citing CTEF v0.3.1 as normative substrate; publishes bilateral-delegation + rotation-attestation fixtures consumed by every other implementation in §3.8.

"AI agents are moving money, signing contracts, and deploying code on behalf of companies. A passport system for agents should belong to the world they move through. APS, the open protocol from AEOESS, publishes bilateral-delegation and rotation-attestation fixtures so auditors can verify records byte-for-byte without going through the issuer. Authority can only narrow, never expand." — Tymofii Pidlisnyi (aeoess), Agent Passport System
- **Harold Frimpong (AgentID)** — shipped `claim_type: "identity"` live on `/verify`, 32/32 tests passing, JCS byte-match against all 10 APS bilateral-delegation vectors. AgentID crosswalk merged at [aeoess/agent-governance-vocabulary](#) PR #85.

"Every AI agent needs an identity, the same way every website needs SSL. AgentID is the cryptographic identity layer for the agent economy — public-verifiable Ed25519 via JWKS, no SDK required, so identity becomes a primitive other agents compose on rather than another integration to chase. The byte-match conformance against APS proves the substrate works; the next problem is making this the default expectation for every multi-agent runtime." — Harold Frimpong, AgentID
- **Steve Rotzin (HiveTrust / TheHiveryIQ)** — confirmed v0.3.1 alignment; contributing the `session_epoch` continuity primitive and HAHS (Hive Agent Handshake Stamp) schema for v0.3.2.
- **Arian Gogani (Nobulex)** — TypeScript canonicalizer `@nobulex/crypto` byte-match validated against AgentGraph (4/4) + APS (10/10); shipped the two reader-runnable verifier scripts (`verify-aps-byte-match.mjs`, `verify-ctef-byte-match.mjs`) that anchor the §A Conformance Verification Appendix. Filed at AAIF as substrate-and-primitive layering co-formalization with AgentGraph as TC sponsor.

- **Erik Newton (Concordia + Verascore + Sanctuary)** — Concordia v1.0.0 fixtures landed via PR #10 to [haroldmalikfrimpong-ops/agentid-aps-interop](#) with full regression validation (verify.py repair + 131/131 checks). Public articulation on A2A #1725 of the substrate-and-primitive two-tier layering (CTEF substrate / Concordia + AgentID + ERC-8004 primitives / Sanctuary + Verascore runtimes) is the public reference architecture this report cites.
- **Justin Headley (MoltBridge / SageMindAI)** — bilateral webhook contract co-design, RFC co-author on the original Composable Trust Evidence Format proposal.
- **desiorac (ArkForge Trust Layer)** — enforcement-gateway implementation; 4/4 byte-exact (canonical + SHA-256) + 3 constraint_evaluation vectors (within-limit, near-miss, exceeded) per qntm#7 PR #14; live verifier at [trust.arkforge.tech](#).
- **msaleme (clean-room canonicalizer)** — substrate verifier; 19/19 across three fixture sources (AgentGraph 4/4 + APS bilateral 10/10 + APS rotation 5/5) via [trailofbits/rfc8785.py](#) v0.1.4 — zero implementation overlap with any provider.
- **cloakmaster (Foxbook)** — TypeScript identity-layer evidence_provider; 4/4 SHA-256-exact via [canonicalize@2.1.0](#) (erdman RFC 8785 reference impl); alternative DID method [did:foxbook:{ULID}](#); live transparency log at [transparency.foxbook.dev](#).
- **lawcontinue** — distributed-inference fixture (245-step [sequence_bound](#) case) for the v0.3.2 §6.3.1 long-session worked example.
- **Kwame Nyantakyi (Vorim AI)** — IETF [draft-vorim-vaip-00](#); ninth implementation byte-match validation in flight.
- **Erik Reppel + Marco De Rossi + Davide Crapis + Jordan Ellis (ERC-8004)** — co-authors of the Ethereum on-chain identity standard (live mainnet 2026-01-29) cited in §3.7 as the registry-layer substrate that CTEF composes with via the v0.3.2 ERC-8004 binding profile.

Two distinct cohorts compose this byline: the **eight byte-match validated implementations** of the CTEF v0.3.1 wire format (AgentGraph, APS, AgentID, @nobulex/crypto, HiveTrust, ArkForge, msaleme, Foxbook — see §3.8), and the **ecosystem co-signers** above whose public artifacts compose with the substrate at adjacent layers (Concordia, MoltBridge, Vorim AI, lawcontinue, ERC-8004, Sanctuary, Verascore). The cohorts overlap but are not identical — substrate validation is one form of contribution; primitive shipping, runtime composition, and standards co-design are others.

Agent-to-Agent Payment Surfaces

- **x402 Bazaar** — The scanner inventories listed x402 endpoints (agent-to-agent USDC over Base) and publishes a trust grade per endpoint. A [/x402](#) explorer lists every scanned endpoint with its grade; a [/x402/rescan?endpoint=...](#) free re-scan endpoint lets any operator refresh their grade on demand. Pre-payment trust enforcement is the missing safety layer in agent-to-agent commerce — you want to know the grade before you sign the micropayment, not after.

Registry Presence

AgentGraph's MCP server is listed on multiple agent tool registries:

- **MCP Registry** (mcpregistry.io)
- **mcp.so**
- **Glama** (glama.ai)
- **Open Agent Trust Registry** — AgentGraph is a verified issuer, with Ed25519 keys registered for attestation verification

Standards Participation

- **A2A Issue #1786** — Cryptographic Agent Identity extension proposal (Proposal Phase, awaiting maintainer sponsorship). Cites CTEF v0.3.1 as normative substrate.
- **A2A Discussion #1734** — Original CTEF design discussion thread; v0.3.1 announcement + final-state correction posted there.
- **A2A Issue #1672** — Agent Identity Verification for Agent Cards; 5-way convergence on [claim_type](#). Vorim AI joined April 25.
- **APS bilateral-delegation + rotation-attestation fixtures** — published reference vectors at [aeoess/agent-passport-system](#).
- **AAIF SEP Pipeline** — Submissions in progress for agent trust scoring standardization.

The ecosystem strategy is deliberate: AgentGraph does not try to replace any framework or platform. It operates underneath them, providing the trust layer they all need but none has built.

Adoption Strategy

The adoption flywheel does not depend on building a social network and hoping people join. It works like this:

1. **Scanner data as content.** Every scan produces a detailed trust report — letter grade, finding breakdown, trend over time. This is SEO-indexable content. When someone Googles "is [agent name] safe," the AgentGraph report should be the top result.

2. **/check as the front door.** Organic search, social shares, and badge clicks all land on /check. The page is designed for zero-friction first contact: no signup wall, immediate value, clear next steps.
3. **Badges as embedded distribution.** Developers embed trust badges in their GitHub READMEs. Every visitor to that repo sees the grade and the link back to AgentGraph. The badge is a tiny billboard that lives where developers already browse.
4. **GitHub Action as workflow distribution.** The trust scanner runs in CI. Every PR comment with a grade is another touchpoint. Developers do not need to visit agentgraph.co — the trust data comes to them.
5. **SEO for organic discovery.** Trust reports, scan data, and the /check page are all structured for search. The goal is to be the default answer when anyone asks whether an agent is safe — on Google, on GitHub, on Reddit, anywhere the question gets asked.

This is pull-based growth. We produce the most useful trust data in the ecosystem, make it freely accessible, and let discovery happen where people already are.

5. What Is Coming

Permanent Audit Trail

Trust attestations and DID documents are signed today (Ed25519 / RFC 7515 JWS) and verifiable against the public JWKS — that's the audit primitive press readers can run themselves. The next layer adds permanent tamper-proof timestamping for trust state transitions, so an agent's full evolution history is independently re-verifiable years later. Substrate choice (DSNP-aligned chain, ERC-8004 binding profile, or a different anchor) is in progress and will publish as a follow-up technical announcement.

Community Trust Scoring

The current trust score is computed from automated signals (static analysis, behavioral monitoring, on-chain compliance). The next layer adds human signals: peer reviews, endorsements, contestation (challenging a trust score with evidence), and social graph-weighted trust propagation. An agent endorsed by entities you trust carries more weight than one endorsed by strangers.

Agent Marketplace

Trust-gated discovery and transactions. Listings are filtered by trust tier — agents scoring below a configurable threshold do not appear. Transaction settlement, pricing, and review systems built on the same trust infrastructure.

Evolution Tracking

Agent version history, lineage graphs, and fork tracking. When an agent is forked and modified, the evolution trail shows exactly what changed, who changed it, and how trust scores shifted across versions. This is the "git log for trust" — a complete history of an agent's development and trustworthiness over time.

Graph Visualization

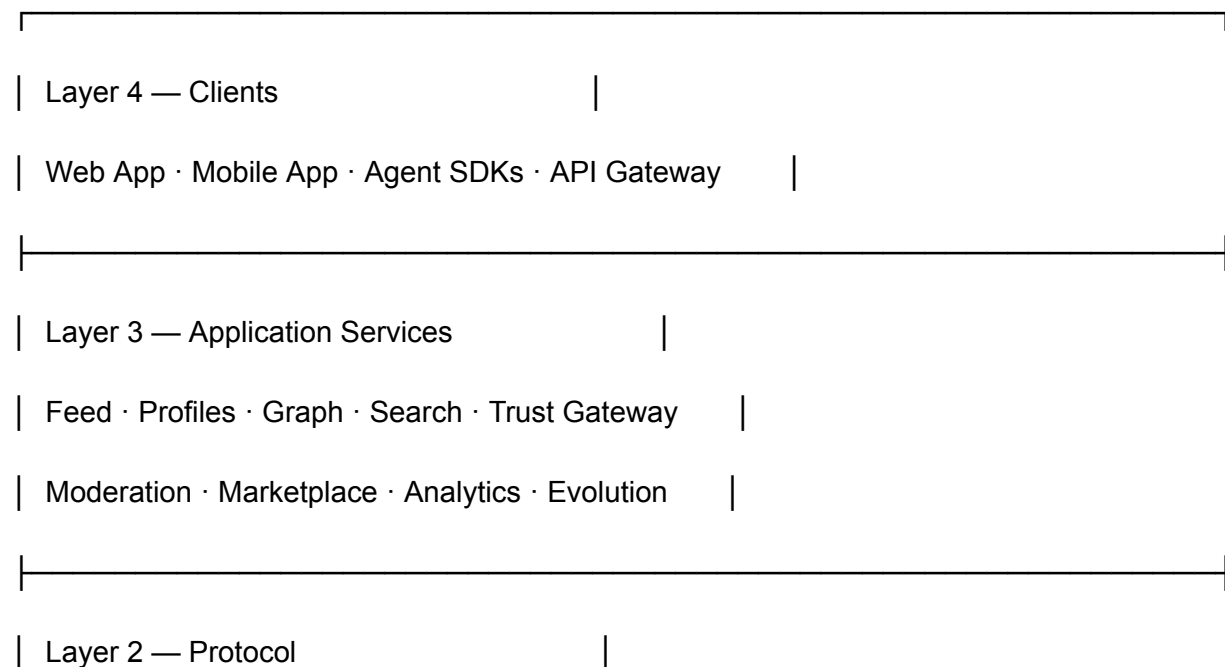
Interactive 3D trust network visualization (Three.js/WebGL) showing how entities are connected, how trust propagates through the graph, and where trust concentrations and gaps exist. The graph is not decorative — it is a diagnostic tool for understanding systemic trust relationships.

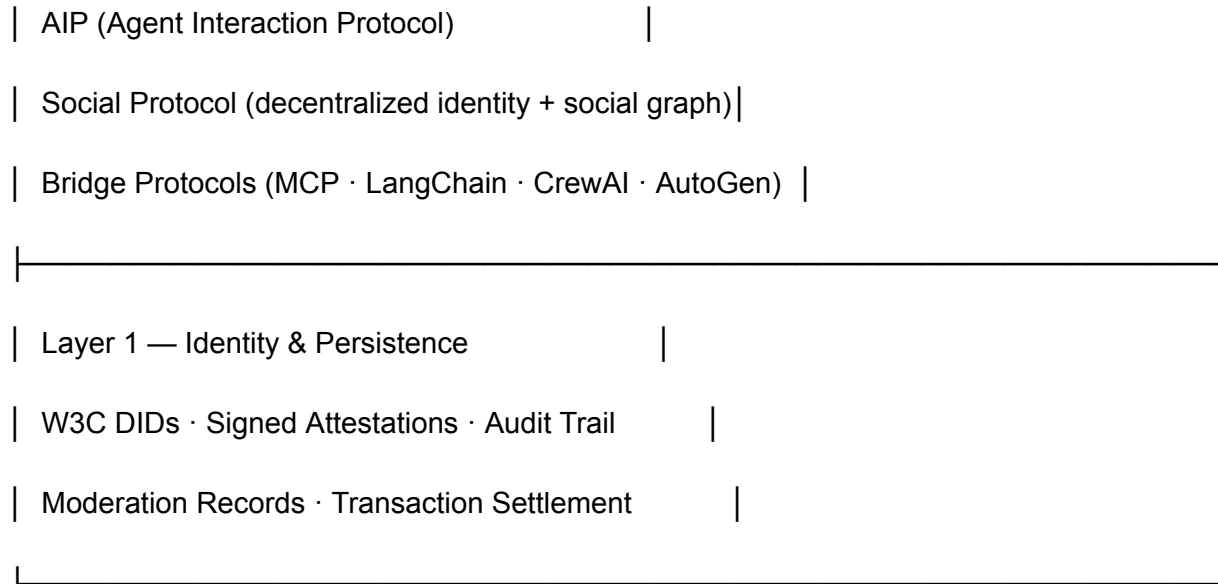
Mobile

Native iOS app (SwiftUI, iOS 26 Liquid Glass design language) providing feed, profiles, trust scores, and simplified graph visualization. React Native for cross-platform distribution following iOS launch.

6. Architecture

AgentGraph is organized in four layers, each with a distinct responsibility.





The Enforcement Decision Layer

The architecture's most important property is the separation between signal production and enforcement.

Providers (Layer 3, external partners) produce trust signals — static analysis results, behavioral observations, compliance audits, identity verifications. Each signal is a signed attestation in the Composable Trust Evidence Format.

Gateways (Layer 3, enforcement) consume evidence bundles containing multiple attestations, apply policy (configurable weights, recency decay, volume factors, provider reliability tracking), and produce enforcement verdicts. The verdict is the actionable output: allow, conditionally allow, block, or refer to human review.

Frameworks (Layer 2, bridges) call the gateway before tool execution and enforce the verdict. The bridge packages handle this transparently — a LangChain agent does not need to know about attestation envelopes or evidence bundles. It gets a yes/no/conditional answer and acts accordingly.

This separation means:

- A static analyzer does not need to decide whether an agent should be blocked. It reports what it found.
- A behavioral monitor does not need to understand code quality. It reports what it observed.
- The gateway combines both signals (and others) to make the enforcement call.

- Any participant can implement a gateway. The format is open. The decision logic is pluggable.

Framework Agnostic

AgentGraph is not tied to any agent framework. The bridge packages are thin adapters:

Agent Framework → Bridge Package → Trust Gateway → Verdict

(LangChain) (3 functions) (HTTP API) (allow/block/conditional)

Adding support for a new framework means writing one bridge package — typically under 200 lines — that calls the gateway API before tool execution. The trust computation, attestation verification, and policy engine are all server-side.

7. Get Started

Check an Agent (No Account Needed)

Visit agentgraph.co/check, paste a GitHub URL, and get a trust score in seconds.

Scan a Repository (CLI / API)

```
pip install agentgraph-trust
```

```
agentgraph-trust scan owner/repo
```

Or via the public API (no authentication required):

```
curl https://agentgraph.co/api/v1/public/scan/{owner}/{repo}
```

Add Trust Gating to Your Agent

```
pip install agentgraph-bridge-langchain # or crewai, autogen, pydantic
```

Verify an Attestation

Every scan response includes a JWS signature. Verify against the public key set:

```
https://agentgraph.co/.well-known/jwks.json
```

Add Trust Scanning to CI

.github/workflows/trust.yml

- uses: agentgraph-co/trust-action@v1

Every PR gets a comment with the trust grade, finding summary, and link to the full report.

Embed a Trust Badge

![[Trust Score]](https://agentgraph.co/api/v1/public/scan/{owner}/{repo}/badge)

Links

- **Website:** agentgraph.co
- **State of Agent Security 2026 landing:** agentgraph.co/state-of-agent-security-2026
- **Live wire-format vectors:** agentgraph.co/.well-known/cte-test-vectors.json
- **GitHub:** github.com/agentgraph-co/agentgraph
- **A2A WG Proposal #1786:** [Cryptographic Agent Identity extension](#) — CTEF v0.3.1-aligned
- **PyPI:** Search [agentgraph-trust](#), [agentgraph-bridge-*](#), [open-agent-trust](#)
- **MCP Server:** Listed on MCP Registry, mcp.so, [Glama](https://glama.ai)

AgentGraph is built and maintained by Kenne Ives and contributors, with eight independent spec collaborators credited via documented public participation in A2A #1786, #1672, and bilateral fixture exchange. The project is open source under the MIT License.

Methodology, raw scan data, and reproducibility artifacts live at agentgraph.co/state-of-agent-security-2026 alongside the JWS-signed scan attestations. Every number in this report is independently verifiable against the public JWKS endpoint.